# CREDENCE

## Augmenting Datacenter Switch Buffer Sharing with ML Predictions

<u>Vamsi Addanki</u>, Maciej Pacut, Stefan Schmid

Technische Universität Berlin

European Research Council
Established by the European Commission

# Let's Play a Game

Event venue

# Let's Play a Game

Entry



A  B  C  D

Event venue

# Let's Play a Game

# Let's Play a Game

Entry

·A   ·B   ·C   ·D

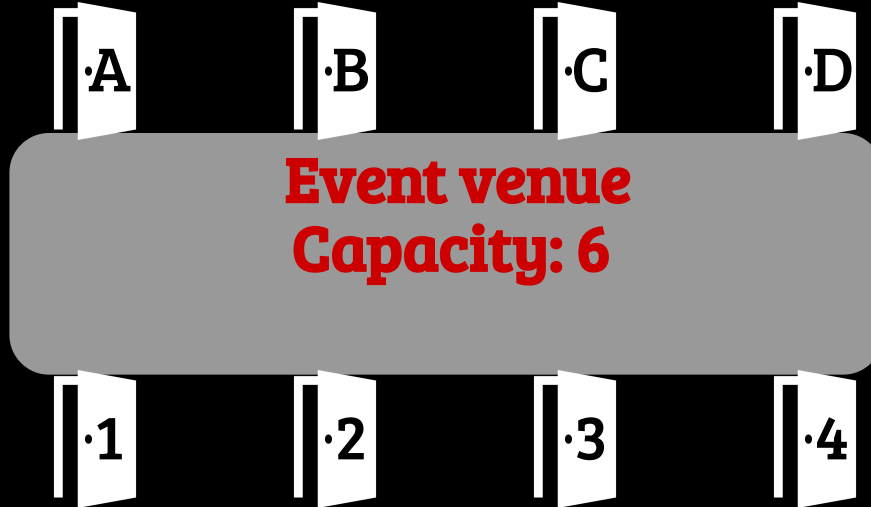**Event venue
Capacity: 6**

Exit

·1   ·2   ·3   ·4

CREDENCE

# Let's Play a Game

Entry

Bouncer

Exit

·A  ·B  ·C  ·D

Event venue
Capacity: 6

·1  ·2  ·3  ·4

CREDENCE

Let's Play a Game

Arrivals: 4

Capacity: 6

A · B · C · D

· 1 · 2 · 3 · 4

CREDENCE

8

# Let's Play a Game

Arrivals: 4

A  B  C  D

Score: 1

Capacity: 6

1  2  3  4

CREDENCE

Let's Play a Game

Arrivals: 8

Score: 1

Capacity: 6

Let's Play a Game

Arrivals: 8

Score: 1

Capacity: 6

A   B   C   D

1   2   3   4

CREDENCE

# Let's Play a Game

Arrivals: 8

A  B  C  D

Score: 2

Capacity: 6

1  2  3  4

CREDENCE

# Let's Play a Game

Arrivals: 12

A  B  C  D

Score: 2

Capacity: 6

1  2  3  4

# Let's Play a Game

Arrivals: 12

A    B    C    D

Score: 2

Capacity: 6

1    2    3    4

CREDENCE

14

# Let's Play a Game

**Arrivals: 16**

**Score: 3**

**Capacity: 6**

A  B  C  D

1  2  3  4

Let's Play a Game

Arrivals: 20

Score: 4

Capacity: 6

A B C D

1 2 3 4

CREDENCE

# Let's Play a Game

**Arrivals: 20**

**Score: 5**

**Capacity: 6**

A    B    C    D

1    2    3    4

CREDENCE

# Let's Play a Game



Previous Game

Arrivals: 20
Score: 5

Capacity: 6

CREDENCE

# Let's Play a Game



Arrivals: 4

A  B  C  D

1  2  3  4

**Previous Game**
Arrivals: 20
Score: 5

Capacity: 6

CREDENCE

# Let's Play a Game

Arrivals: 4

A    B    C    D

Score: 1

1    2    3    4

Capacity: 6

CREDENCE

# Let's Play a Game

Arrivals: 12

Score: 3

·A   ·B   ·C   ·D

·1   ·2   ·3   ·4

Capacity: 6

CREDENCE

# Let's Play a Game

Arrivals: 12

·A  ·B  ·C  ·D

Score: 3

·1  ·2  ·3  ·4

Capacity: 6

**Previous Game**
Arrivals: 20
Score: 5

CREDENCE

# Let's Play a Game

Arrivals: 16

A  B  C  D

Score: 7

Capacity: 6

1  2  3  4

**Previous Game**
Arrivals: 20
Score: 5

CREDENCE

32

# Let's Play a Game

**Arrivals: 16**

**Score: 7**

·A  ·B  ·C  ·D

·1  ·2  ·3  ·4

**Previous Game**
**Arrivals: 20**
**Score: 5**

**Capacity: 6**

CREDENCE

33

# Let's Play a Game

Arrivals: 20

Score: 11

A B C D

1 2 3 4

**Previous Game**
Arrivals: 20
Score: 5

Capacity: 6

CREDENCE

35

# Let's Play a Game

Arrivals: 20

Score: 11

·A  ·B  ·C  ·D

·1  ·2  ·3  ·4

Capacity: 6

CREDENCE

# Let's Play a Game

**New Game**
Arrivals: 20
Score: 15

·A  ·B  ·C  ·D

·1  ·2  ·3  ·4

**Previous Game**
Arrivals: 20
Score: 5

CREDENCE

# Let's Play a Game

**New Game**
**Arrivals: 20**
**Score: 15**

·A    ·B    ·C    ·D

**Admission Control Algorithms
can improve throughput
*(or severely impact throughput)***

·1    ·2    ·3    ·4

**Previous Game**
**Arrivals: 20**
**Score: 5**

CREDENCE

Buffer Sharing

Input Ports

A    B    C    D

Shared Buffer

Buffer Sharing Algorithm

Network Switch

1    2    3    4

Output Ports

CREDENCE

40

# Buffer Sharing: An Emerging Critical Problem

- Bursty traffic requires buffers to avoid packet losses
- Stringent performance requirements
- But buffer sizes are unable to scale with capacity increase

CREDENCE

# Buffer Sharing: An Emerging Critical Problem

- Bursty traffic requires buffers to avoid packet losses
- Stringent performance requirements
- But buffer sizes are unable to scale with capacity increase

Buffer Sharing algorithm can severely impact end-to-end performance e.g., FCTs

CREDENCE

# Buffer Sharing (An Online Perspective)

- **Goal:** Maximize the number of transmitted packets
  - Throughput maximization

CREDENCE

# Buffer Sharing (An Online Perspective)

- **Goal:** Maximize the number of transmitted packets
  - Throughput maximization
- **Online algorithm (ALG)** takes spontaneous decisions upon every packet arrival

CREDENCE

# Buffer Sharing (An Online Perspective)

- **Goal:** Maximize the number of transmitted packets
  - Throughput maximization
- **Online algorithm (ALG)** takes spontaneous decisions upon every packet arrival
- **Offline optimal algorithm (OPT)** has prior knowledge of the entire arrival sequence and performs optimally

CREDENCE

# Buffer Sharing (An Online Perspective)

- ALG is $C$ competitive if OPT transmits no more than $C$ times that of ALG
  - *OPT ≤ C · ALG*

# Buffer Sharing (An Online Perspective)

- ALG is $C$ competitive if OPT transmits no more than $C$ times that of ALG
  - $OPT \leq C \ ALG$

**Competitive Ratio**

# Online Buffer Sharing Algorithms

- **Drop-tail:** Drop on arrival or accept
  - All commodity switches support drop-tail buffers
- **Push-out:** Accept all packets and push a packet out when the buffer is full
  - *Not supported in hardware*

CREDENCE

Competitive ratio

1 → N

Optimal Throughput — Lower Throughput

Harmonic

Dynamic Thresholds

Complete Sharing

1    Competitive ratio    N

Optimal Throughput

Lower Throughput

CREDENCE

Longest Queue Drop (Push-out)

Harmonic

Dynamic Thresholds

Complete Sharing

1     1.707     **Competitive ratio**     N

Optimal Throughput

Lower Throughput

CREDENCE

Not supported in hardware!

All commodity switches support drop-tail

LQD (Push-out)

Harmonic

Dynamic Thresholds

Complete Sharing

1          1.707          Competitive ratio          N

Optimal Throughput

Lower Throughput

CREDENCE

Not supported in hardware!

All commodity switches support drop-tail

Can we unlock this space?

LQD (Push-out)

Harmonic

Dynamic Thresholds

Complete Sharing

1    1.707    Competitive ratio    N

Optimal Throughput

Lower Throughput

CREDENCE

53

# Limitations of Traditional Drop-tail Buffer Sharing

- **Proactive unnecessary drops** → **throughput loss**

CREDENCE

# Limitations of Traditional Drop-tail Buffer Sharing

- **Proactive unnecessary drops** → <span style="color:red">**throughput loss**</span>
  - Overprovisioning for the *unknown* future arrivals
  - Packet drops are unnecessary if the future is known

CREDENCE

# Limitations of Traditional Drop-tail Buffer Sharing

- **Proactive unnecessary drops** →**throughput loss**
  - Overprovisioning for the *unknown* future arrivals
  - Packet drops are unnecessary if the future is known

- **Reactive avoidable drops** →**throughput loss**

# Limitations of Traditional Drop-tail Buffer Sharing

- **Proactive unnecessary drops** ⟶ <span style="color:red">throughput loss</span>
  - Overprovisioning for the *unknown* future arrivals
  - Packet drops are unnecessary if the future is known

- **Reactive avoidable drops** ⟶ <span style="color:red">throughput loss</span>
  - Underprovisioning for the unknown future arrivals
  - Packet drops are avoidable if the future is known

CREDENCE

**Push-out** **Drop-tail**

LQD
(Push-out)

Harmonic

Dynamic
Thresholds

Complete
Sharing

1          1.707          Competitive ratio          N

Optimal
Throughput

Lower
Throughput

CREDENCE

# Predictions: A Hope for Competitive Buffer Sharing

- **Predict the actions of a push-out algorithm (LQD)**
- **Augment drop-tail algorithms with predictions**
  - **Peek into the future**

CREDENCE

# Predictions: A Hope for Competitive Buffer Sharing

- **Predict the actions of a push-out algorithm (LQD)**
- **Augment drop-tail algorithms with predictions**
  - **Peek into the future**

**Can predictions improve drop-tail's competitive ratio?**

# Naive Approach

- ## Upon a packet arrival
  - ### Predict LQD's action
  - ### If prediction is to accept, then accept
  - ### If prediction is to drop, then drop

CREDENCE

# Challenge: Imperfect Predictions

### True Positive
**Ground Truth: Drop**
**Prediction: Drop**

### False Negative
**Ground Truth: Drop**
**Prediction: Accept**

### False Positive
**Ground Truth: Accept**
**Prediction: Drop**

### True Negative
**Ground Truth: Accept**
**Prediction: Accept**

CREDENCE

# Challenge: Imperfect Predictions

- **Excessive false positives can lead to starvation**
  - *eg., every prediction is "drop"*

CREDENCE

# Challenge: Imperfect Predictions

- **Excessive false positives can lead to starvation**
  - *eg., every prediction is "drop"*
- **Even a single false negative can hurt throughput forever**
  - **(discussed in the paper)**

CREDENCE

# Goals

- **Consistency (under perfect predictions)**
  - Competitive ratio close to push-out
- **Robustness (with large prediction error)**
  - Competitive ratio close to existing algorithms
- **Smoothness**
  - Competitive ratio smoothly degrades with prediction error

CREDENCE

Not supported in hardware!

All commodity switches
support drop-tail

Can we unlock this space?

LQD
(Push-out)

Harmonic          DT          CS

1          1.707          Competitive ratio          N

Optimal
Throughput

Lower
Throughput

CREDENCE

*without predictions*

Push-out ← → Drop-tail

LQD
(Push-out)

Harmonic     DT     CS

1     1.707     Competitive ratio     N

Optimal
Throughput

Lower
Throughput

CREDENCE

Perfect predictions ← **CREDENCE** → Large prediction error

*Drop-tail with predictions*

- - - - - - - - - - - - - - - - - - - - - - - - - - - -

*without predictions*

← Push-out | Drop-tail →

LQD (Push-out)

Harmonic    DT    CS

1    1.707    Competitive ratio    N

Optimal Throughput

Lower Throughput

CREDENCE

68

# Credence

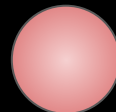- Drop-tail buffer sharing augmented with predictions
- Threshold-based (similar to existing algorithms)

# Credence

- **Drop-tail buffer sharing augmented with predictions**
- **Threshold-based (similar to existing algorithms)**
- **Consistency** ✅
  - **Close to push-out under perfect predictions**

CREDENCE

# Credence

- **Drop-tail buffer sharing augmented with predictions**
- **Threshold-based (similar to existing algorithms)**
- **Consistency** ✅
  - **Close to push-out under perfect predictions**
- **Robustness** ✅
  - **Close to existing algorithms even with large prediction error**

CREDENCE

# Credence

- **Drop-tail buffer sharing augmented with predictions**
- **Threshold-based (similar to existing algorithms)**
- **Consistency ✅**
  - **Close to push-out under perfect predictions**
- **Robustness ✅**
  - **Close to existing algorithms even with large prediction error**
- **Smoothness ✅**
  - **Smoothly degrades with prediction error**

CREDENCE

# Credence's Thresholds

- **Per-queue thresholds**
  - Thresholds are incremented and decremented based on Longest Queue Drop (Push-out) algorithm

CREDENCE

# Credence's Thresholds

- **Per-queue thresholds**
  - Thresholds are incremented and decremented based on Longest Queue Drop (Push-out) algorithm
- **A packet is rejected immediately if the queue length is greater than its corresponding threshold**

CREDENCE

# Credence's Thresholds

- **Per-queue thresholds**
  - Thresholds are incremented and decremented based on Longest Queue Drop (Push-out) algorithm
- **A packet is rejected immediately if the queue length is greater than its corresponding threshold**
- **A <span style="color:red">prediction</span> is obtained *only if* the queue length is lower than its corresponding threshold**

CREDENCE

# Credence's Thresholds

- **Thresholds enable tackling false negative errors**
  - **Prevents accepting too many packets** eg., if all the predictions are "accept"

CREDENCE

# Credence's Thresholds

- **Thresholds enable tackling false negative errors**
  - **Prevents accepting too many packets** eg., if all the predictions are "accept"
- **Safe guard criterion to tackle false positive errors**
  - Always accept a packet if the longest queue is lower than fair-share of buffer partition
  - **Prevents dropping too many packets** eg., if all the predictions are "drop"
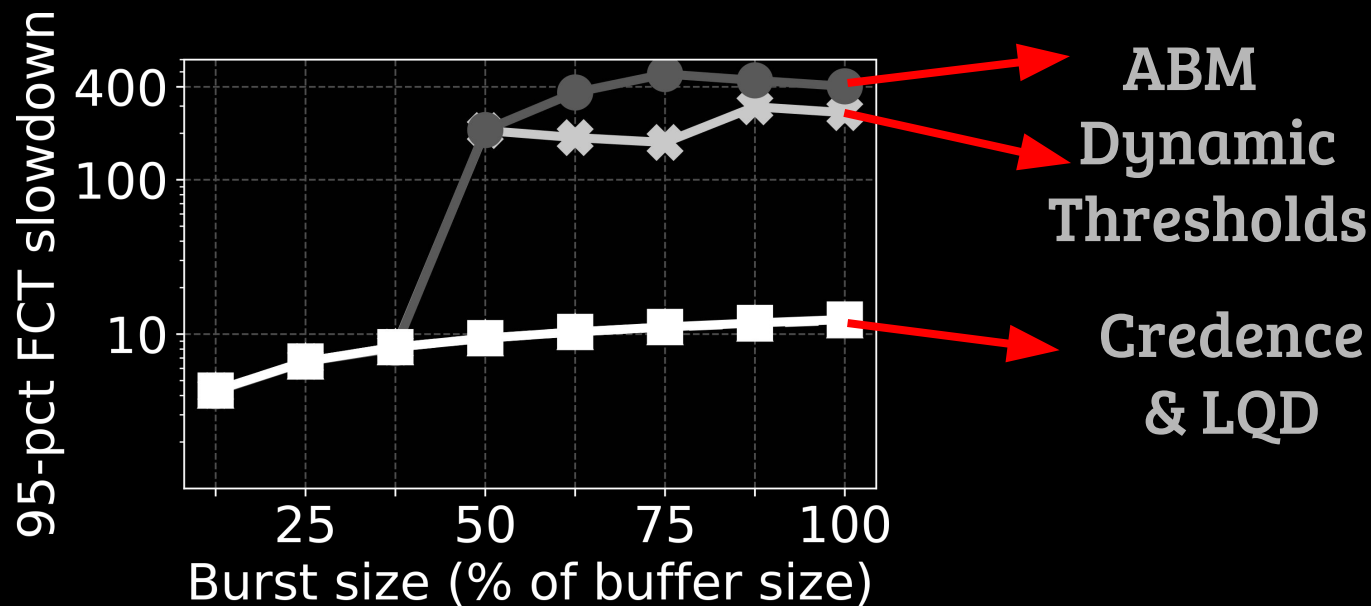
CREDENCE

# Further Details in the Paper

- Competitive analysis
- Theoretical bounds for Credence's performance
- ...

CREDENCE

# Evaluation

- Packet-level simulations using NS3
- 256 servers, 4 spine switches and 16 ToR switches
- 10Gbps NICs
- Shared buffer at the switches
- Random Forest-based prediction oracle for Credence

CREDENCE

# Credence Performs Close to Push-out



Chart: 95-pct FCT slowdown (y-axis, log scale with marks at 10, 100, 400) vs Burst size (% of buffer size) (x-axis, marks at 25, 50, 75, 100)

- ABM
- Dynamic Thresholds
- Credence & LQD

CREDENCE

# Credence Degrades with Prediction Error

# Open Questions and Future Research Directions

- **Practically training a prediction oracle**
  - Simulation-based data (may not capture real-world scenarios)
  - Real-world network data (more accurate but complex to obtain)
    - Online reinforcement learning
    - …
- **Understanding push-out operation complexity**
- **Improving the robustness of Credence**
- **Considering latency for competitive analysis**

CREDENCE

# Conclusion

- Traditional drop-tail buffer sharing approaches cannot be improved further
- Credence is the first buffer sharing algorithm augmented with predictions
- Credence offers bounded performance guarantees
- Credence can improve the performance of datacenter traffic in terms of flow completion times for short flows and incast flows
- Source code: https://github.com/inet-tub/ns3-datacenter

CREDENCE

**Vamsi Addanki**
*vamsi@inet.tu-berlin.de*
🐦 @Vamsi_DT

**Maciej Pacut**
*maciej@inet.tu-berlin.de*

**Stefan Schmid**
*stefan.schmid@tu-berlin.de*
🐦 @schmiste_ch

# Thank You